

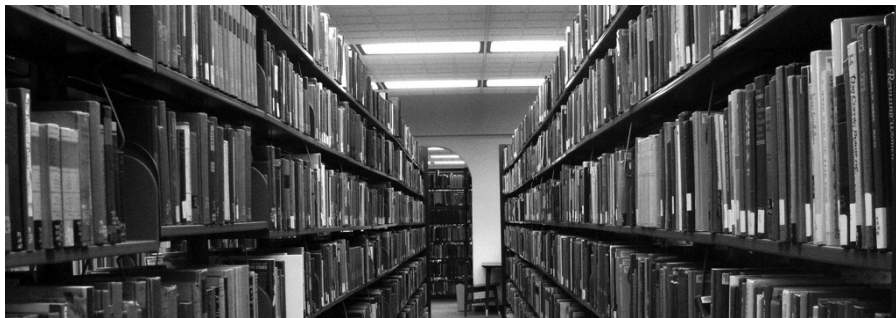
# **Probabilistic Topic Models and User Behavior**

David M. Blei

Columbia University



- ▶ **ORGANIZE**
- ▶ **VISUALIZE**
- ▶ **SUMMARIZE**
- ▶ **SEARCH**
- ▶ **PREDICT**
- ▶ **UNDERSTAND**



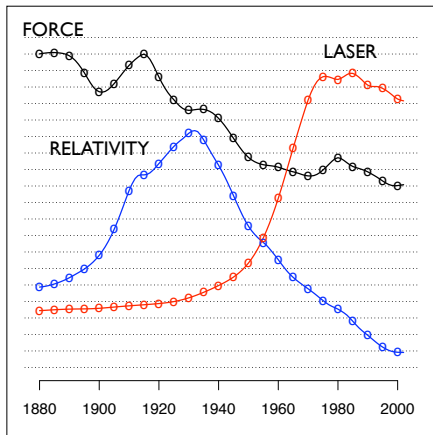
## TOPIC MODELING

1. **Discover** the thematic structure
2. **Annotate** the documents
3. **Use** the annotations to visualize, organize, summarize, ...

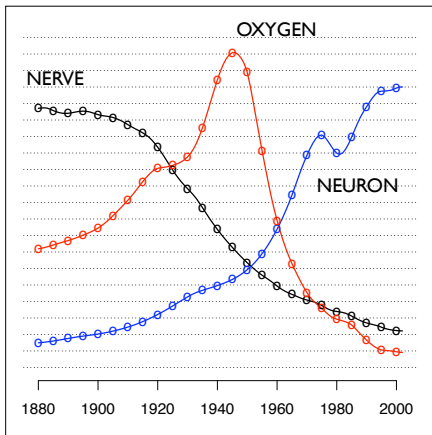


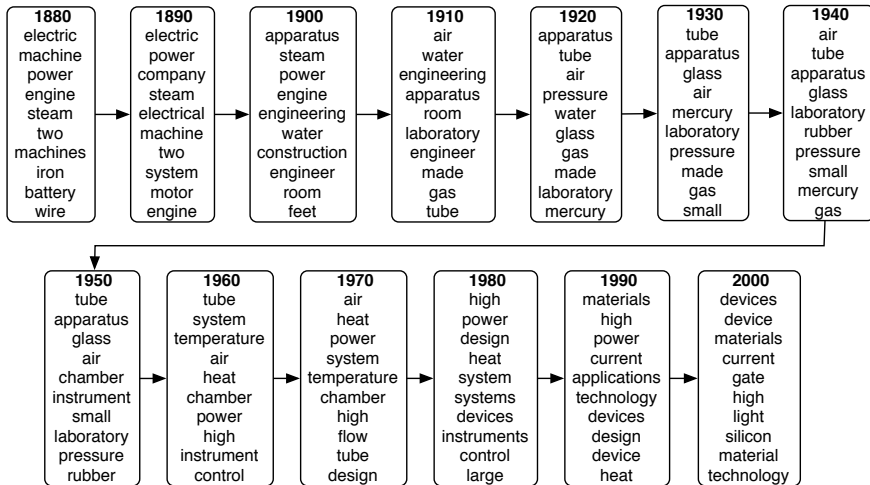


## "Theoretical Physics"



## "Neuroscience"







SKY WATER TREE  
MOUNTAIN PEOPLE



SCOTLAND WATER  
FLOWER HILLS TREE



SKY WATER BUILDING  
PEOPLE WATER



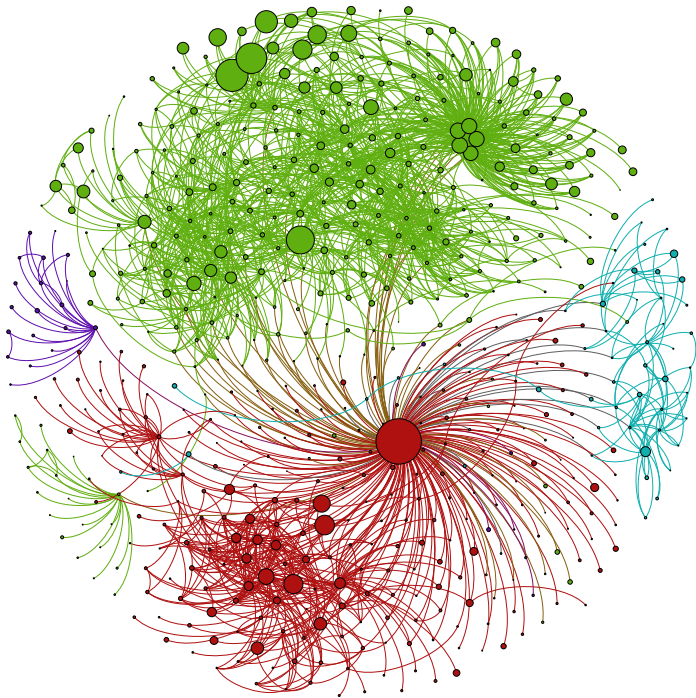
FISH WATER OCEAN  
TREE CORAL

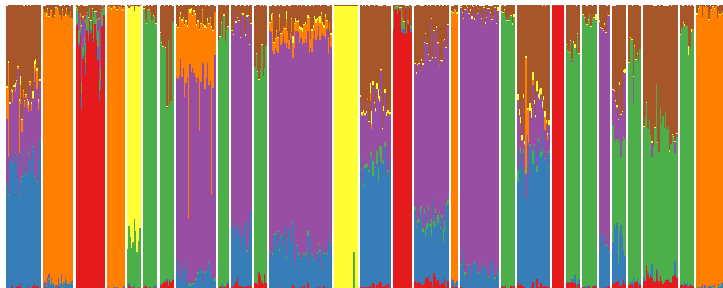
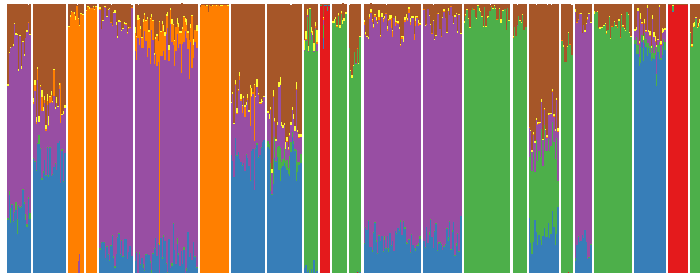


PEOPLE MARKET PATTERN  
TEXTILE DISPLAY



BIRDS NEST TREE  
BRANCH LEAVES







Charles Darwin's library



The NYC subway

- ▶ **People read documents.**
- ▶ These might be people for whom we want to form predictions.
- ▶ And, their behavior is an additional signal about the meaning of the documents and the organization of the collection.

## **This talk**

1. Introduction to topic modeling
2. Recommendation and exploration with collaborative topic models
3. The bigger picture: Using probability models to solve problems with data

# **Introduction to Topic Modeling**



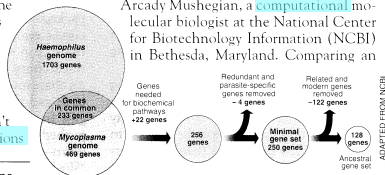
## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers game**, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains

Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Documents exhibit multiple topics.

## Topics

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

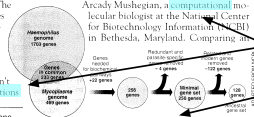
## Documents

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here, two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson, a geneticist at the University of Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly if more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

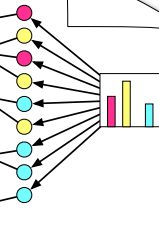


\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

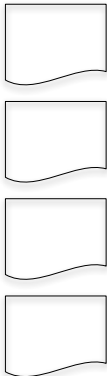
SCIENCE • VOL. 272 • 24 MAY 1996

## Topic proportions and assignments



## Latent Dirichlet Allocation

## Topics



## Documents

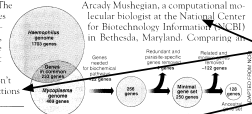
## Topic proportions and assignments

### Seeking Life's Bare (Genetic) Necessities

**COLD SPRING HARBOR, NEW YORK**—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

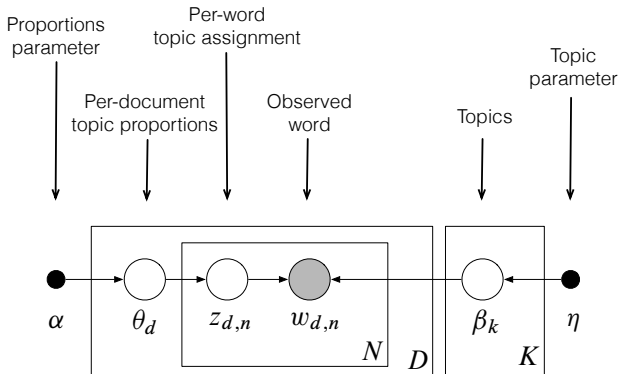
"are not all that far apart," especially in comparison to the 75,000 genes to the human genome, notes Siv Anderson, a biologist at the University of Warwick, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing the



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

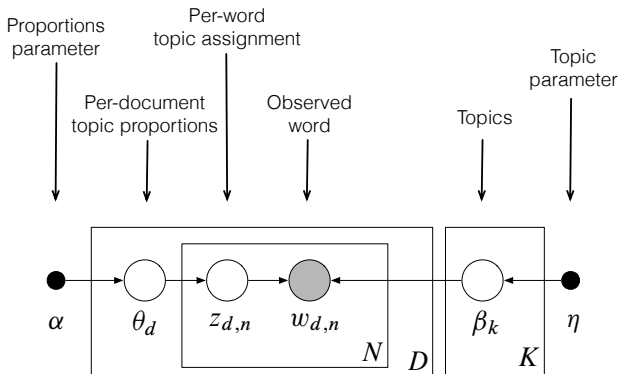
SCIENCE • VOL. 272 • 24 MAY 1996

## Latent Dirichlet Allocation



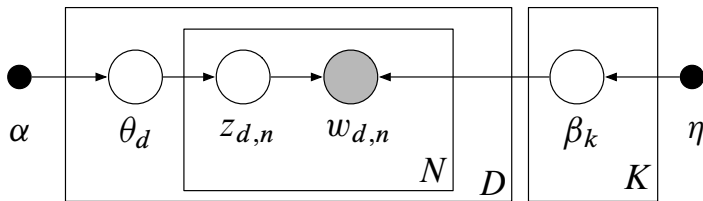
### LDA as a graphical model

- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; unshaded nodes are hidden.
- Plates indicate replicated variables.

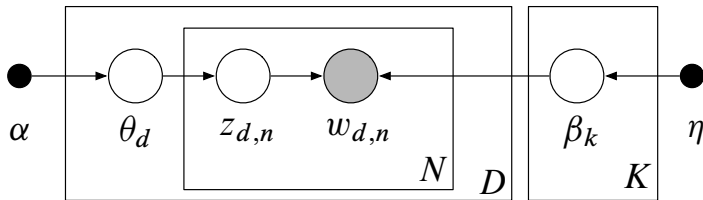


### LDA as a graphical model

- Encodes independence assumptions about the variables
- Defines a factorization of the joint probability distribution
- Connects to algorithms for computing with data



- The joint defines a posterior,  $p(\theta, z, \beta \mid w)$ .
- From a collection of documents, infer
  - Per-word topic assignment  $z_{d,n}$
  - Per-document topic proportions  $\theta_d$
  - Per-corpus topic distributions  $\beta_k$
- Then use posterior expectations to perform the task at hand: information retrieval, document similarity, exploration, and others.



- ▶ Mean field variational methods (Blei et al., 2001, 2003)
- ▶ Expectation propagation (Minka and Lafferty, 2002)
- ▶ Collapsed Gibbs sampling (Griffiths and Steyvers, 2002)
- ▶ Distributed sampling (Newman et al., 2008; Ahmed et al., 2012)
- ▶ Collapsed variational inference (Teh et al., 2006)
- ▶ Stochastic inference (Hoffman et al., 2010, 2013; Mimno et al., 2012)
- ▶ Factorization inference (Arora et al., 2012; Anandkumar et al., 2012)



- ▶ **Data:** The OCR'ed collection of *Science* from 1990–2000
  - 17K documents
  - 11M words
  - 20K unique terms (stop words and rare words removed)
- ▶ **Model:** 100-topic LDA model using variational inference.

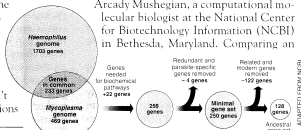


## Seeking Life's Bare (Genetic) Necessities

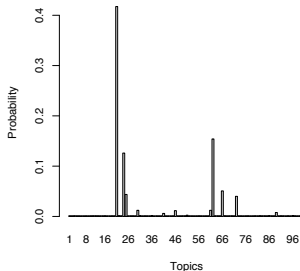
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.



\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

1

Game  
Season  
Team  
Coach  
Play  
Points  
Games  
Giants  
Second  
Players

2

Life  
Know  
School  
Street  
Man  
Family  
Says  
House  
Children  
Night

3

Film  
Movie  
Show  
Life  
Television  
Films  
Director  
Man  
Story  
Says

4

Book  
Life  
Books  
Novel  
Story  
Man  
Author  
House  
War  
Children

5

Wine  
Street  
Hotel  
House  
Room  
Night  
Place  
Restaurant  
Park  
Garden

6

Bush  
Campaign  
Clinton  
Republican  
House  
Party  
Democratic  
Political  
Democrats  
Senator

7

Building  
Street  
Square  
Housing  
House  
Buildings  
Development  
Space  
Percent  
Real

8

Won  
Team  
Second  
Race  
Round  
Cup  
Open  
Game  
Play  
Win

9

Yankees  
Game  
Mets  
Season  
Run  
League  
Baseball  
Team  
Games  
Hit

10

Government  
War  
Military  
Officials  
Iraq  
Forces  
Iraqi  
Army  
Troops  
Soldiers

11

Children  
School  
Women  
Family  
Parents  
Child  
Life  
Says  
Help  
Mother

12

Stock  
Percent  
Companies  
Fund  
Market  
Bank  
Investors  
Funds  
Financial  
Business

13

Church  
War  
Women  
Life  
Black  
Political  
Catholic  
Government  
Jewish  
Pope

14

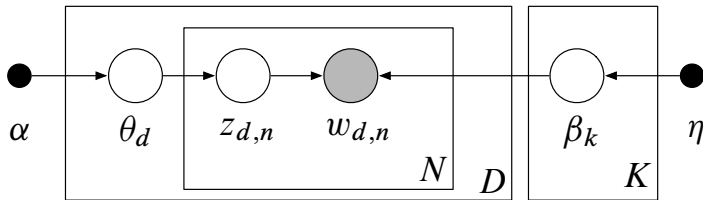
Art  
Museum  
Show  
Gallery  
Works  
Artists  
Street  
Artist  
Paintings  
Exhibition

15

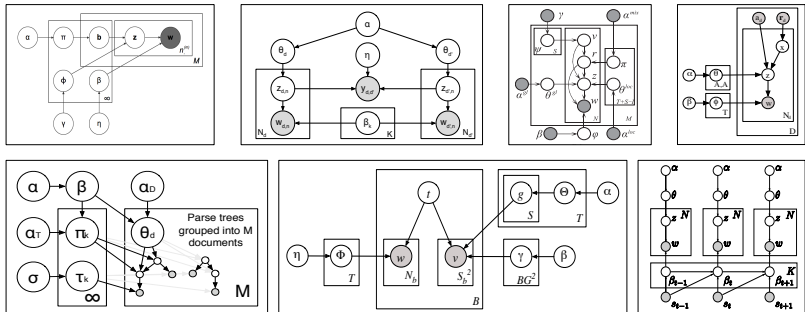
Police  
Yesterday  
Man  
Officer  
Officers  
Case  
Found  
Charged  
Street  
Shot

## How does LDA “work”?

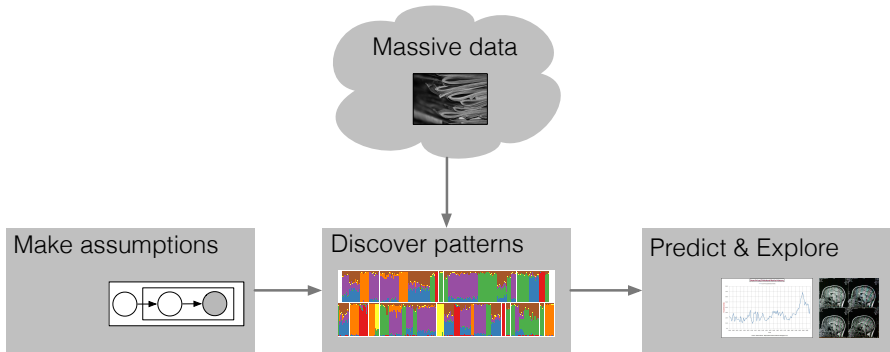
- ▶ LDA trades off two goals.
  1. In each **document**, allocate its words to **few topics**.
  2. In each **topic**, assign high probability to **few terms**.
- ▶ These goals are at odds.
  - Putting a document in a single topic makes #2 hard:  
All of its words must have probability under that topic.
  - Putting very few words in each topic makes #1 hard:  
To cover a document's words, it must assign many topics to it.
- ▶ Trading off these goals finds groups of tightly co-occurring words.



- LDA discovers themes through posterior inference.
- Other perspectives
  - Latent semantic analysis [Deerwester et al., 1990; Hofmann, 1999]
  - A mixed-membership model [Erosheva, 2004]
  - PCA and matrix factorization [Jakulin and Buntine, 2002]
  - Was independently invented for genetics [Pritchard et al., 2000]



- Topic modeling is an active field of research.  
LDA is a simple building block that enables many applications.
- Organizing and finding patterns in text has become important in the sciences, humanities, industry, and culture.
- Algorithmic improvements let us fit models to massive data.  
(See VW, Gensim, Mallet, others.)



- ▶ Case study in **text analysis with probability models**
- ▶ Topic modeling research
  - develops new models.
  - develops new inference algorithms.
  - develops new applications, visualizations, tools.

**Quick diversion: Fitting a topic model in your spare time**



perspective identifying tumor suppressor genes in human...  
letters global warming report leslie roberts article global....  
research news a small revolution gets under way the 1990s....  
a continuing series the reign of trial and error draws to a close...  
making deep earthquakes in the laboratory lab experimenters...  
quick fix for freeways thanks to a team of fast working...  
feathers fly in grouse population dispute researchers...

....

245 1897:1 1467:1 1351:1 731:2 800:5 682:1 315:6 3668:1 14:1  
260 4261:2 518:1 271:6 2734:1 2662:1 2432:1 683:2 1631:7  
279 2724:1 107:3 518:1 141:3 3208:1 32:1 2444:1 182:1 250:1  
266 2552:1 1993:1 116:1 539:1 1630:1 855:1 1422:1 182:3 2432:1  
233 1372:1 1351:1 261:1 501:1 1938:1 32:1 14:1 4067:1 98:2  
148 4384:1 1339:1 32:1 4107:1 2300:1 229:1 529:1 521:1 2231:1  
193 569:1 3617:1 3781:2 14:1 98:1 3596:1 3037:1 1482:12 665:2

....

```
docs <- read.documents("mult.dat")  
K <- 20  
alpha <- 1/20  
eta <- 0.001  
model <- lda.collapsed.gibbs.sampler(documents, K, vocab, 1000, alpha, eta)
```

## LDA in R

1 dna gene sequence genes sequences human genome genetic analysis two	2 protein cell cells proteins receptor fig binding activity activation kinase	3 water climate atmospheric temperature global surface ocean carbon atmosphere changes	4 says researchers new university just science like work first years	5 mantle high earth pressure seismic crust temperature earths lower earthquakes
6 end article start science readers service news card circle letters	7 time data two model fig system number element reads see	8 materials surface high structure temperature molecules chemical molecular fig university	9 dna rna transcription protein site binding sequence proteins specific sequences	10 disease cancer patients human gene medical studies drug normal drugs
11 years million ago age university north early fig evidence record	12 species evolution population evolutionary university populations natural studies genetic biology	13 protein structure proteins two amino binding acid residues molecular structural	14 cells cell virus hiv infection immune human antigen infected viral	15 space solar observations earth stars university mass sun astronomers telescope
16 fax manager science aas advertising sales member recruitment associate washington	17 cells cell gene genes expression development mutant mice fig biology	18 energy electron state light quantum physics electrons high laser magnetic	19 research science national scientific scientists new states university united health	20 neurons brain cells activity fig channels university cortex neuronal visual

## Wikipedia Topics

Relative Presence of Topics in all Documents

{household, population, female}

{film, series, show}

{theory, work, human}

{son, year, death}

{war, force, army}

{system, computer, user}

{album, band, music}

{government, party, election}

{game, team, player}

{god, call, give}

{company, market, business}

{math, number, function}

{few, some, great}

## {film, series, show}

words	related documents	related topics
film	The X-Files	{son, year, death}
series	Orson Welles	{work, book, publish}
show	Stanley Kubrick	{album, band, music}
character	B movie	{woman, child, man}
play	Mystery Science Theater 3000	{law, state, case}
make	Monty Python	{black, white, people}
episode	Doctor Who	{theory, work, human}
movie	Sam Peckinpah	{{@card@}, make, design}
good	Married... with Children	{war, force, army}
release	History of film	{god, call, give}
feature	The A-Team	{game, team, player}
television	Pulp Fiction (film)	{day, year, event}
star	Mad (magazine)	{company, market, business}

## Stanley Kubrick



### related topics

{film, series, show}  
 {theory, work, human}  
 {son, year, death}  
 {black, white, people}  
 {god, call, give}  
 {math, energy, light}

**Stanley Kubrick** (July 26, 1928 – March 7, 1999) was an American film director, writer, producer, and photographer who lived in England during most of the last four decades of his career. Kubrick was noted for the scrupulous care with which he chose his subjects, his slow method of working, the variety of genres he worked in, his technical perfectionism, and his reluctance about his films and personal life. He worked far beyond the confines of the Hollywood system, maintaining almost complete artistic control and making movies according to his own whims and time constraints, but with the rare advantage of big-studio financial support for all his endeavors.

Kubrick's films are characterized by a formal visual style and meticulous attention to detail—his later films often have elements of surrealism and expressionism that eschews structured linear narrative. His films are repeatedly described as slow and methodical, and are often perceived as a reflection of his obsessive and perfectionist nature.<sup>[1]</sup> A recurring theme in his films is man's inhumanity to man. While often viewed as

### related documents

Orson Welles  
 B movie  
 Mystery Science Theater 3000  
 Monty Python  
 Doctor Who  
 Sam Peckinpah  
 The A-Team  
 Pulp Fiction (film)  
 Buffy the Vampire Slayer (TV series)  
 The X-Files  
 Sunset Boulevard (film)  
 Jack Benny

## {theory, work, human}

words	related documents	related topics
theory	Meme	{work, book, publish}
work	Intelligent design	{law, state, case}
human	Immanuel Kant	{son, year, death}
idea	Philosophy of mathematics	{woman, child, man}
term	History of science	{god, call, give}
study	Free will	{black, white, people}
view	Truth	{film, series, show}
science	Psychoanalysis	{war, force, army}
concept	Charles Peirce	{language, word, form}
form	Existentialism	{{@card@}, make, design}
world	Deconstruction	{church, century, christian}
argue	Social sciences	{rate, high, increase}
social	Idealism	{company, market, business}

# **Collaborative Topic Models**

with Prem Gopalan, Laurent Charlin, and Chong Wang



Charles Darwin's library



Reading on the New York subway

- ▶ **People read documents.**
- ▶ *Collaborative topic models* connect content to consumption

# Exploration and Exploitation of Victorian Science in Darwin's Reading Notebooks

Jaimie Murdock<sup>1,2</sup>, Colin Allen<sup>1,3,4</sup>, and Simon DeDeo<sup>\*1,2,5</sup>

<sup>1</sup>Program in Cognitive Science, Indiana University, Bloomington, IN 47405, USA

<sup>2</sup>School of Informatics and Computing, Indiana University, Bloomington, IN 47405, USA

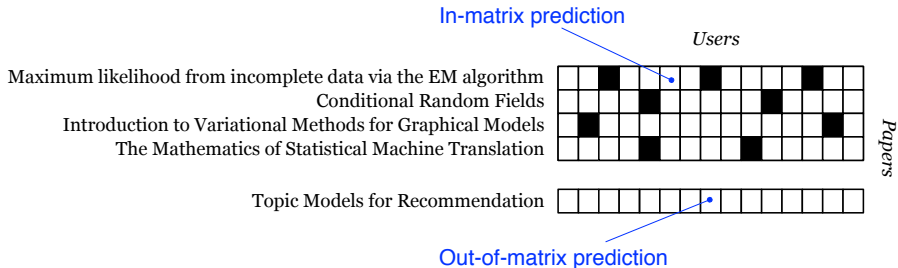
<sup>3</sup>Department of History and Philosophy of Science and Medicine, Indiana University, Bloomington, IN 47405, USA

<sup>4</sup>School of Humanities and Social Sciences, Xi'an Jiaotong University, Xi'an, China

<sup>5</sup>Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

September 25, 2015

<http://arxiv.org/abs/1509.07175>



- ▶ Example: Scientists share their research libraries.
- ▶ Collaborative topic models can
  - Helps readers discover documents
  - Describe readers in terms of topical preferences
  - Identify impactful, interdisciplinary articles

- Consider EM (Dempster et al., 1977). We infer topics from its text:

Maximum Likelihood from Incomplete Data via the *EM* Algorithm

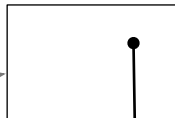
By A. P. DEMPSTER, N. M. LAIRD and D. B. RUBIN

*Harvard University and Educational Testing Service*

[Read before the ROYAL STATISTICAL SOCIETY at a meeting organised by the RESEARCH SECTION on Wednesday, December 8th, 1976, Professor S. D. SILVEY in the Chair]

SUMMARY

A broadly applicable algorithm for computing maximum likelihood estimates from incomplete data is presented at various levels of generality. Theory showing the monotone behaviour of the likelihood and convergence of the algorithm is derived. Many examples are sketched, including missing value situations, applications to grouped, censored or truncated data, finite mixture models, variance component estimation, hyperparameter estimation, iteratively reweighted least squares and factor analysis.



Vision Statistics

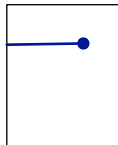
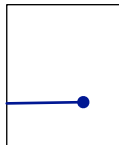
- Suppose there are two types of scientists

STATISTICIAN

VISION RESEARCHER

Vision

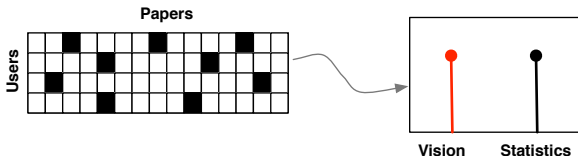
Statistics



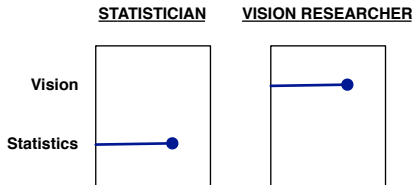
- We first recommend the EM paper to **statisticians**.



- With user data, we can adjust the topics to account for who liked it:



- Consider again the scientists



- We now recommend the EM paper to **vision researchers**.

## Maximum Likelihood from Incomplete Data via the *EM* Algorithm

By A. P. DEMPSTER, N. M. LAIRD and D. B. RUBIN

*Harvard University and Educational Testing Service*

[Read before the ROYAL STATISTICAL SOCIETY at a meeting organized by the RESEARCH SECTION on Wednesday, December 8th, 1976, Professor S. D. SILVEY in the Chair]

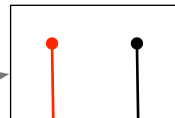
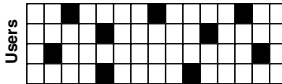
### SUMMARY

A broadly applicable algorithm for computing maximum likelihood estimates from incomplete data is presented at various levels of generality. Theory showing the monotone behaviour of the likelihood and convergence of the algorithm is derived. Many examples are sketched, including missing value situations, applications to grouped, censored or truncated data, finite mixture models, variance component estimation, hyperparameter estimation, iteratively reweighted least squares and factor analysis.



Vision      Statistics

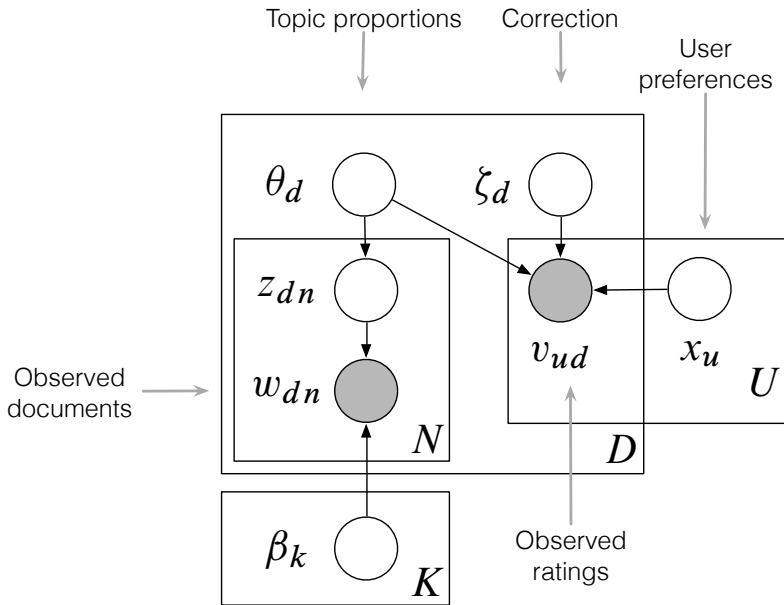
## Papers



Vision      Statistics

**Without text, we cannot initially recommend to anyone.**

**Without user data, we cannot recommend to vision researchers.**



# Maximum Likelihood from Incomplete Data via the EM Algorithm

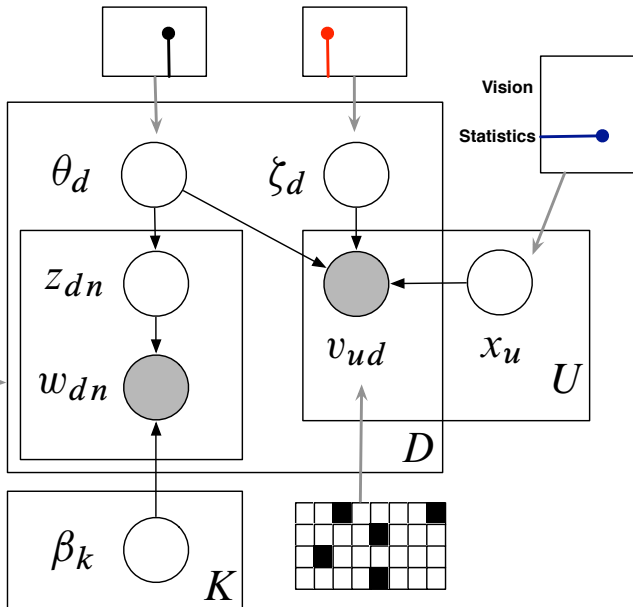
By A. P. Dempster, N. M. Laird and D. B. Rubin

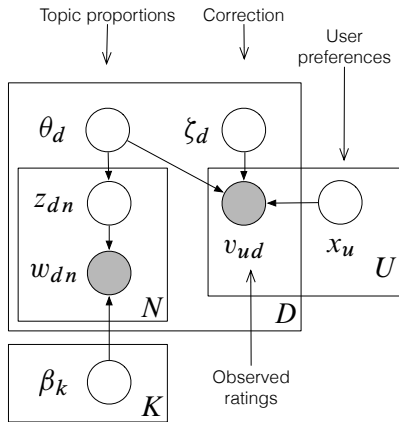
Harvard University and Educational Testing Service

Read before the Royal Statistical Society at a meeting organized by the Research Section on Wednesday, December 10, 1979, Professor J. D. Sacks in the Chair

## Synopsis

A broadly applicable algorithm for computing maximum likelihood estimates from incomplete data is presented at various levels of generality. Theory showing the monotone behaviour of the likelihood and convergence of the algorithm is derived. Many examples are sketched, including missing value situations, applications to grouped, censored or truncated data, finite mixture models, variance component estimation, hyperparameter estimation, iteratively reweighted least squares and factor analysis.





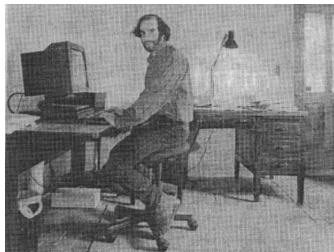
$$\begin{aligned}
 \theta_{dk} &\sim \text{Gamma}(\cdot, \cdot) \\
 \xi_{dk} &\sim \text{Gamma}(\cdot, \cdot) \\
 x_{uk} &\sim \text{Gamma}(\cdot, \cdot) \\
 v_{ud} &\sim \text{Poisson}((\theta_d + \xi_d)^\top x_u)
 \end{aligned}$$

- Blends factorization-based and content-based recommendation
- Describes user preferences with interpretable topics
- Builds on Poisson factorization

[Canney 2004; Dunson and Herring 2005; Gopalan et al. 2014)



- ▶ Big data set from Mendeley.com
- ▶ The data:
  - 261K documents
  - 80K users
  - 10K vocabulary terms
  - 25M observed words
  - 5.1M entries (sparsity is 0.02%)



- ▶ A decade of clicks on arXiv.org (2003–2013)
- ▶ The data:
  - 826K documents
  - 120K users
  - 14K vocabulary terms
  - 54M observed words
  - 43.6M entries (sparsity is 0.04%)

# Maximum Likelihood from Incomplete Data via the *EM* Algorithm

By A. P. DEMPSTER, N. M. LAIRD and D. B. RUBIN

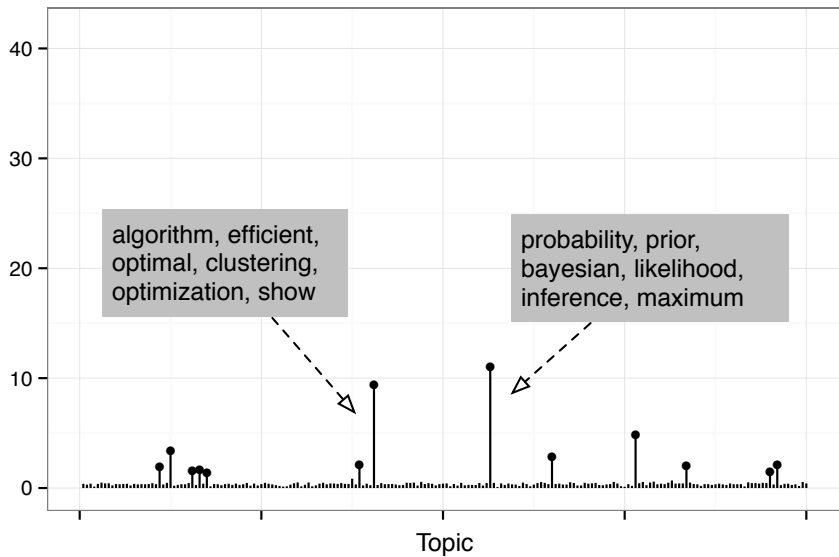
*Harvard University and Educational Testing Service*

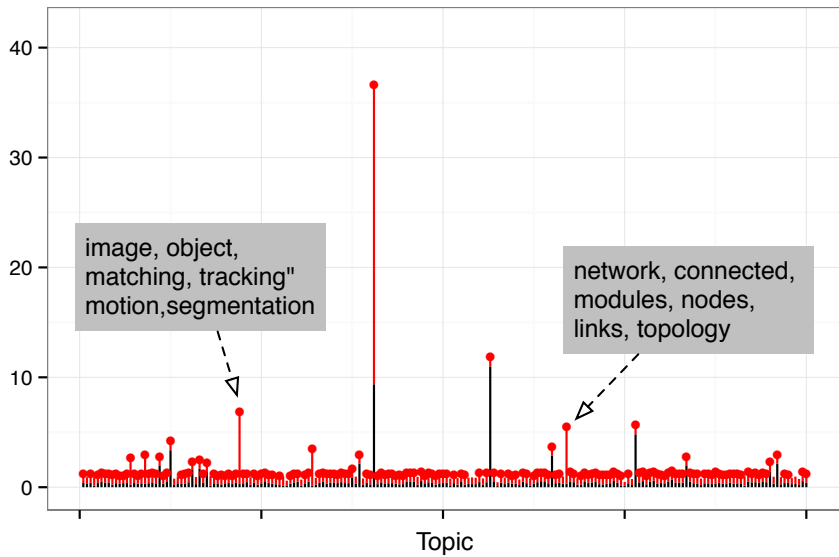
[Read before the ROYAL STATISTICAL SOCIETY at a meeting organized by the RESEARCH SECTION on Wednesday, December 8th, 1976, Professor S. D. SILVEY in the Chair]

## SUMMARY

A broadly applicable algorithm for computing maximum likelihood estimates from incomplete data is presented at various levels of generality. Theory showing the monotone behaviour of the likelihood and convergence of the algorithm is derived. Many examples are sketched, including missing value situations, applications to grouped, censored or truncated data, finite mixture models, variance component estimation, hyperparameter estimation, iteratively reweighted least squares and factor analysis.



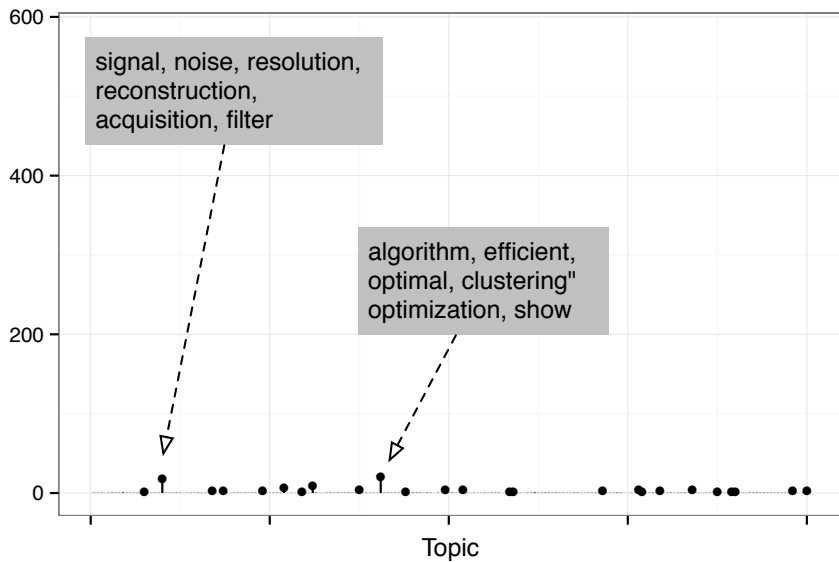


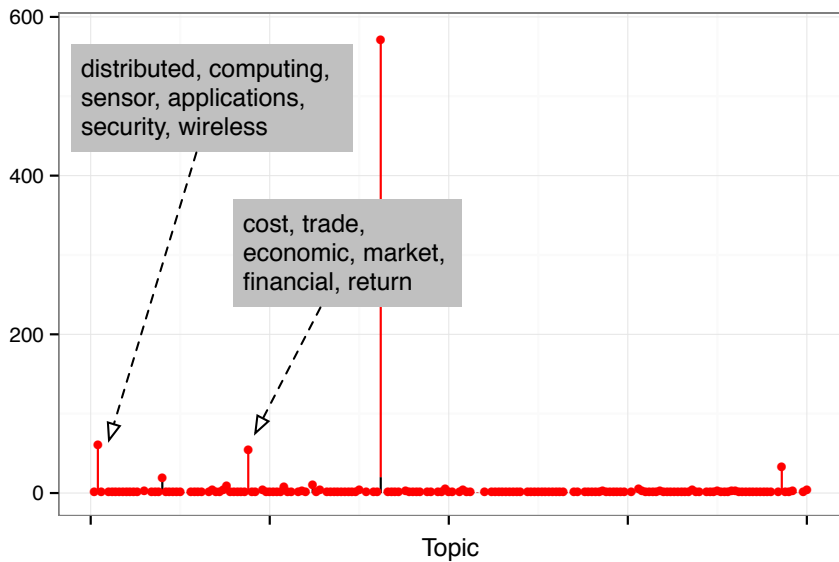


Stephen Boyd and  
Lieven Vandenberghe

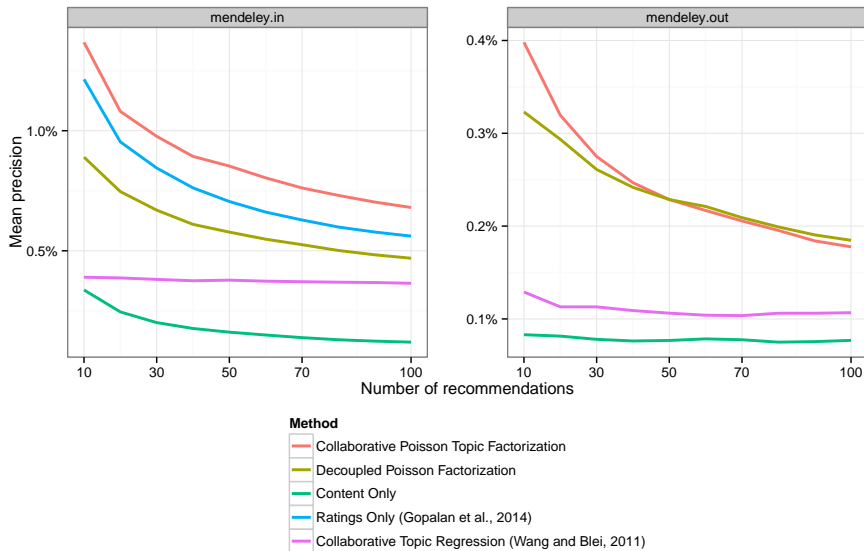
# Convex Optimization

CAMBRIDGE

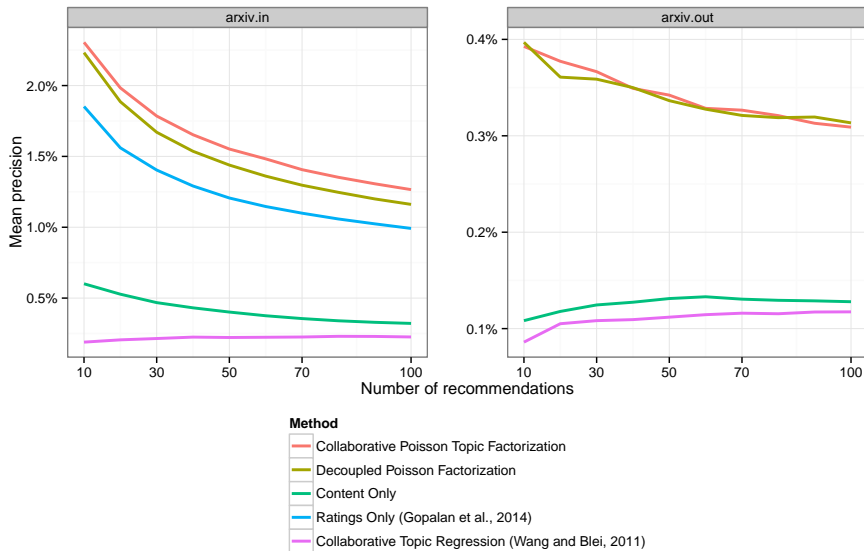




# Mendeley

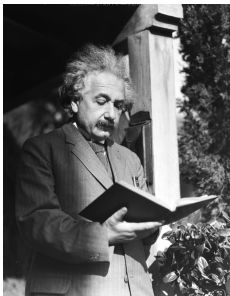


# arXiv click history

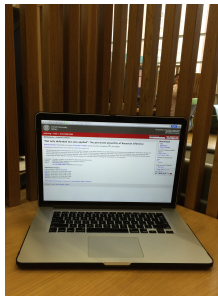




Darwin's library



Einstein reading

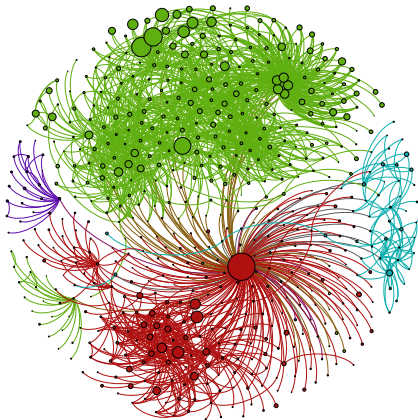


Another scientist reading

- ▶ The readers also **tell us about the articles.**
- ▶ We can look at posterior estimates to find
  - Interdisciplinary articles
  - Influential articles within a field
  - Outside influences on a field



# “Network Analysis”



network; connected; modules; nodes; links; topology; connectivity; graph;  
robustness; connections; modular; world; degree; properties

## Assortative mixing in networks

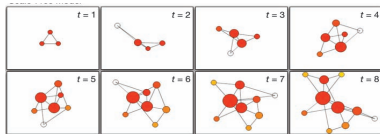
M. E. J. Newman

*Department of Physics, University of Michigan, Ann Arbor, MI 48109-1120 and  
Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501*



## About networks

- ▶ Assortative mixing in networks  
(Newman, 2002)
- ▶ Mixing patterns in networks  
(Newman, 2002)
- ▶ Catastrophic cascade of failures in interdependent networks  
(Buldyrev et al., 2010)



## About networks; for readers of networks

- ▶ Emergence of scaling in random networks  
(Barabassi and Albert, 1999)
- ▶ Statistical mechanics of complex networks  
(Albert and Barabassi, 2002)
- ▶ Complex networks: Structure and dynamics  
(Boccaletti et al., 2006)

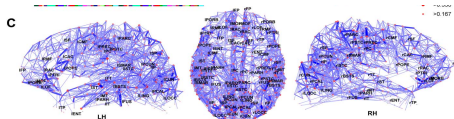
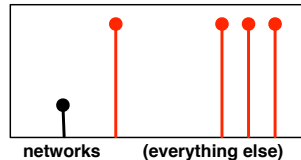
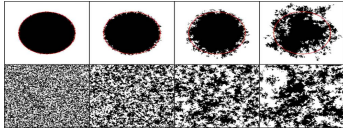


Figure 3. High-Resolution Connection Matrix, Network Layout and Connectivity Backbone (Participant A, scan 2)



## About networks; for readers of other fields

- ▶ Mapping the Structural Core of Human Cerebral Cortex  
(Hagmann et al., 2008)
- ▶ Network thinking in ecology and evolution  
(Proulx et al., 2005)
- ▶ Linked: The New Science of Networks  
(Barabasi, 2002)



## ***Not about networks; for readers of networks***

- ▶ Power-law distributions in empirical data  
(Clauset et al., 2009)
- ▶ Statistical physics of social dynamics  
(Castellano et al., 2009)
- ▶ The origin of bursts and heavy tails in human dynamics  
(Barabasi, 2005)

# “Statistical Modeling”

## **About this field; read by users in this field**

- ▶ A Bayesian analysis of some nonparametric problems
- ▶ Bayesian measures of model complexity and fit
- ▶ Monte Carlo Methods in Bayesian Computation

## **About this field; read by users in other fields**

- ▶ A tutorial on HMMs and selected applications in speech recognition
- ▶ An Introduction to Bayesian Networks and Influence Diagrams
- ▶ Maximum likelihood from incomplete data via the EM algorithm

## **About other fields; read by users in this field**

- ▶ Second Thoughts on the Bootstrap
- ▶ A guide to Eclipse and the R plug-in StatET
- ▶ Using Multivariate Statistics

# “Algorithms”

## **About algorithms; for readers of algorithms**

- ▶ Convex Optimization
- ▶ Nonlinear dimensionality reduction by locally linear embedding
- ▶ Independent component analysis: Algorithms and applications

## **About algorithms; for readers of other fields**

- ▶ Introduction to Algorithms
- ▶ Fast approximate energy minimization via graph cuts
- ▶ Nonlinear dimensionality reduction by locally linear embedding

## **About other fields; for readers of algorithms**

- ▶ A Mathematical Theory of Communication
- ▶ An Introduction To Compressive Sampling
- ▶ Elements of Information Theory

# “Information Theory” (from the ArXiv data)

## **About information theory; for readers of information theory**

- ▶ Wireless Network Information Flow: A Deterministic Approach
- ▶ Random Access: An Information-Theoretic Perspective
- ▶ Coding for Network Coding

## **About information theory; for readers of other fields**

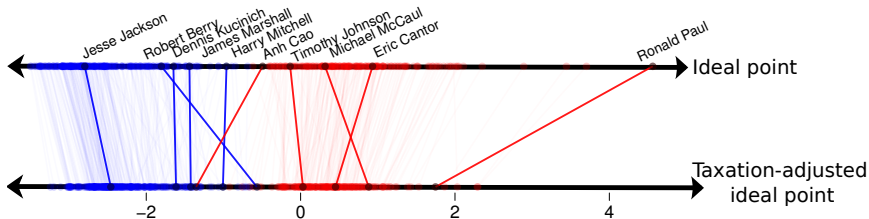
- ▶ Quantum Communication With Zero-Capacity Channels
- ▶ Can non-private channels transmit quantum information?
- ▶ Quantum Steganography

## **About other fields; for readers of information theory**

- ▶ Matrix Completion from a Few Entries
- ▶ Message Passing Algorithms for Compressed Sensing
- ▶ Adaptive Alternating Minimization Algorithms



## The issue-adjusted ideal point model [Gerrish and Blei, 2012]



- ▶ **Roll call data** is also a type of user behavior data.
- ▶ Classical matrix factorization captures the coarse political spectrum.
- ▶ But lawmakers can deviate based on the issues of the bill.
- ▶ The **issue-adjusted ideal point model** captures this deviation.
- ▶ On taxation Ron Paul (R) is more liberal than expected  
Robert Berry (D) is more conservative than expected.



## Collaborative topic models

- ▶ Connect text to usage, content to consumption
- ▶ Blend content-based and user-based recommendation
- ▶ Opens new windows into how people read

## **Discussion: Modern Probabilistic Modeling**

TOPIC  
MODELING

The diagram consists of two nested ellipses. The outer ellipse is white with a black border. Inside it is a smaller, light-gray ellipse with a black border. The text 'STATISTICS', 'MACHINE LEARNING', and 'DATA SCIENCE' is centered at the bottom of the outer ellipse. Inside the gray ellipse, the text 'TOPIC MODELING' is on the left and 'PROBABILISTIC MODELING' is at the bottom. An arrow points from 'TOPIC MODELING' to a small black dot located in the upper right area of the gray ellipse.

PROBABILISTIC  
MODELING

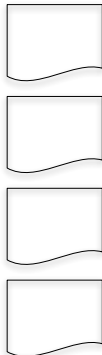
STATISTICS  
MACHINE LEARNING  
DATA SCIENCE

# I. Assume our data come from a model with hidden patterns at work

Topics

Documents

Topic proportions and assignments

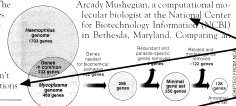


## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson, a biologist at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a numeric game. Some particularly "core" and more genomes are repeatedly targeted and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing the



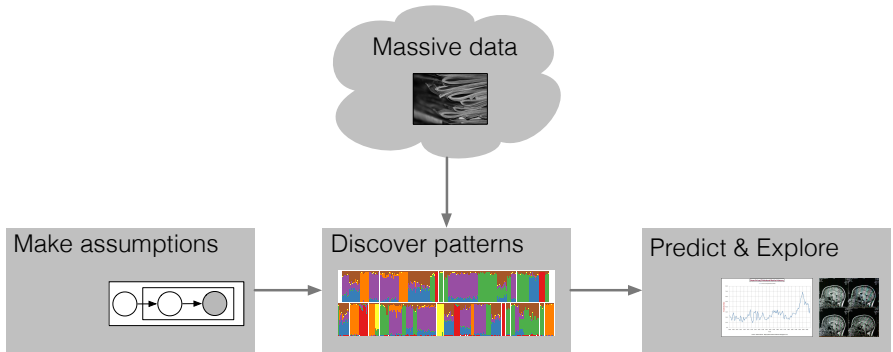
\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

## II. Discover those patterns from data

$$\nu^* = \arg \max_{\nu} \mathbb{E}_q [\log p(x, z, \beta \mid \alpha)] + \mathbb{H} [q(z, \beta \mid \nu)]$$

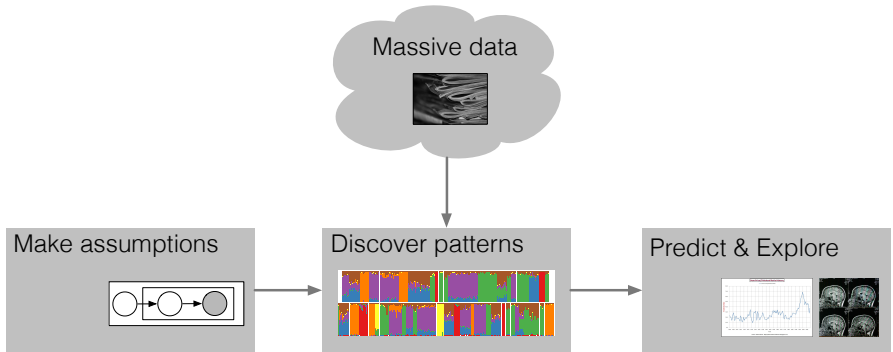




Our perspective:

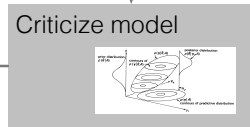
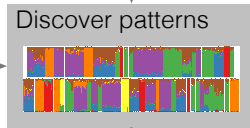
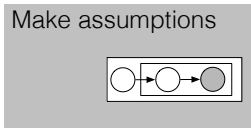
- ▶ Customized data analysis is important to many fields.
- ▶ This pipeline separates assumptions, computation, application.
- ▶ It facilitates solving data science problems.



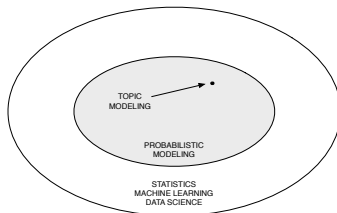
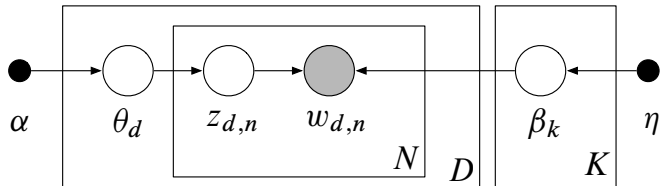


What we need:

- ▶ **Flexible** and **expressive** components for building models
- ▶ **Scalable** and **generic** inference algorithms
- ▶ **Easy to use** software to stretch probabilistic modeling into new areas



Revise

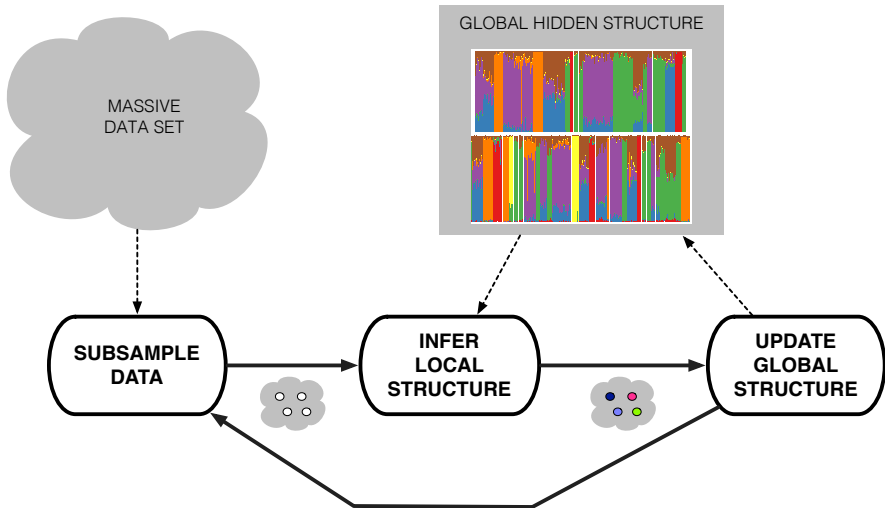




We should seek out unfamiliar summaries of observational material, and establish their useful properties... And still more novelty can come from finding, and evading, still deeper lying constraints.

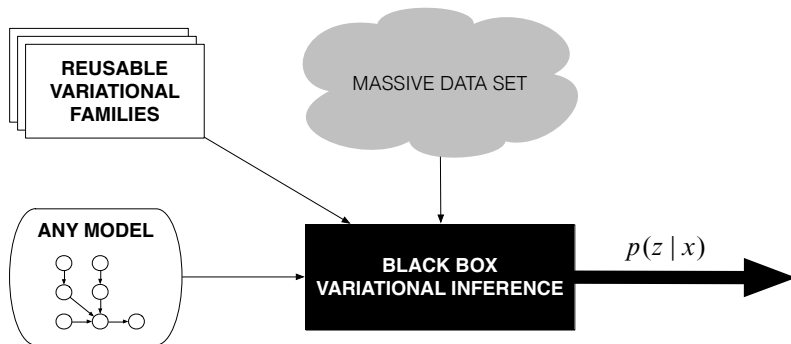
(John Tukey, *The Future of Data Analysis*, 1962)

**A few slides about inference**

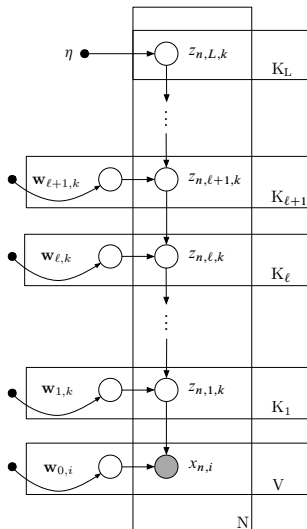


**Stochastic Variational Inference** [Hoffman et al., 2013]

## Black box variational inference



- ▶ Easily use variational inference with *any model*
- ▶ No exponential family requirements
- ▶ No mathematical work beyond specifying the model



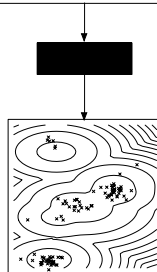
## Deep Exponential Families

[Ranganath et al., 2015]

```
data {
  int<lower=1> K;
  int<lower=1> N;
  real y[N];
}

parameters {
  simplex[K] theta;
  real mu[K];
  real<lower=0,upper=10> sigma[K];
}

model {
  real ps[K];
  for (k in 1:K) {
    mu[k] ~ normal(0,10);
  }
  for (n in 1:N) {
    for (k in 1:K) {
      ps[k] <- log(theta[k])
        + normal_log(y[n],mu[k],sigma[k]);
    }
    lp__ <- lp__ + log_sum_exp(ps);
  }
}
```



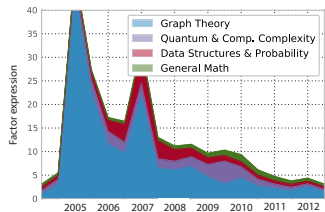
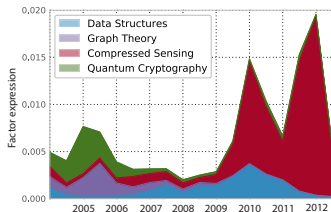
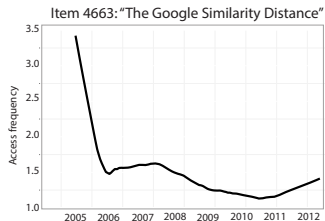
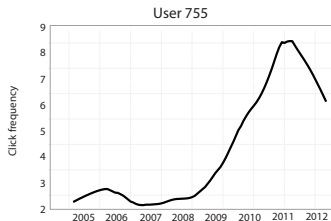
## Probabilistic Programming

[Kucukelbir et al., 2015]



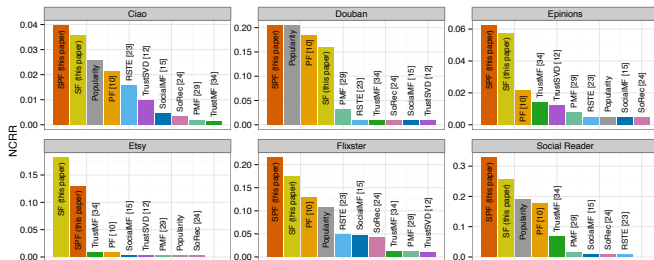
**Some recent work about recommendation**

# Time series recommendation



(with Laurent Charlin, James McInerney, Rajesh Ranganath)

# Social recommendation



(with Allison Chaney)